

Regresión lineal: enfoque Bayesiano

Matías Carrasco

14 de septiembre de 2023

Índice

1. La distribución gaussiana multivariada	1
2. Regresión lineal bayesiana	3
2.1. El modelo	3
2.2. Verosimilitud	4
2.3. A priori	4
2.4. A posteriori	4
2.5. Comparación con la regularización Ridge	6
2.6. Predicciones	6
2.7. Verosimilitud marginal	7

1. La distribución gaussiana multivariada

Un vector de D atributos $\mathbf{x} = [x^{(1)}, \dots, x^{(D)}]^\top$ tiene distribución gaussiana de **media $\boldsymbol{\mu}$** y **matriz de covarianzas $\boldsymbol{\Sigma}$** si su densidad en \mathbb{R}^D está dada por la siguiente fórmula:

$$\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

La media no es otra cosa que el vector de medias $\boldsymbol{\mu} = [\mu^{(1)}, \dots, \mu^{(D)}]^\top$ de cada atributo y la matriz de covarianzas está dada por

$$\boldsymbol{\Sigma} = \mathbb{E} [\mathbf{x}\mathbf{x}^\top] - \boldsymbol{\mu}\boldsymbol{\mu}^\top$$

y es por lo tanto una matriz simétrica y positiva.

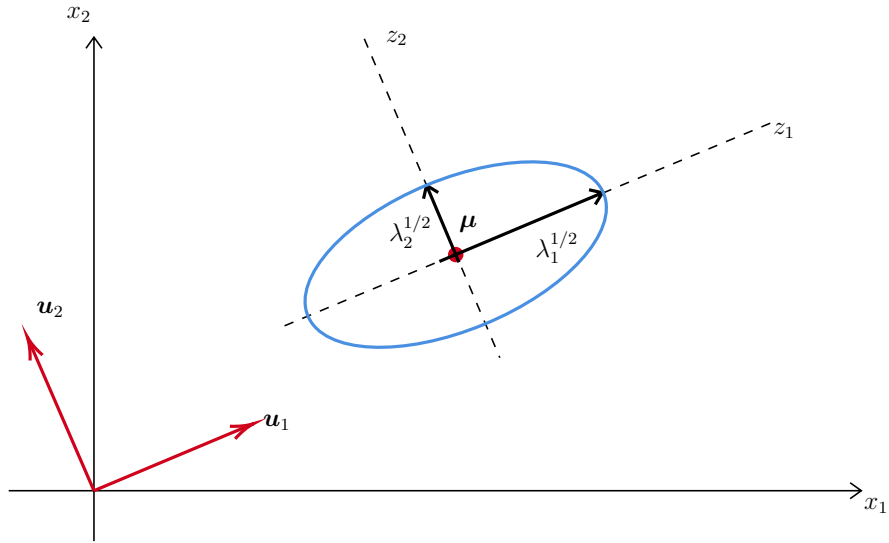


Figura 1: Curvas de nivel de la distribución gaussiana multivariada.

Existe una matriz ortogonal $\mathbf{U} = [\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(D)}]$ y una matriz diagonal $\mathbf{D} = [\lambda^{(1)}, \dots, \lambda^{(D)}]$ tales que

$$\mathbf{\Sigma} = \mathbf{U}\mathbf{D}\mathbf{U}^\top = \sum_{j=1}^D \lambda^{(j)} \mathbf{u}^{(j)} (\mathbf{u}^{(j)})^\top$$

Con esta notación es fácil ver que la inversa es:

$$\mathbf{\Sigma}^{-1} = \mathbf{U}\mathbf{D}^{-1}\mathbf{U}^\top = \sum_{j=1}^D \frac{1}{\lambda^{(j)}} \mathbf{u}^{(j)} (\mathbf{u}^{(j)})^\top$$

Considerando las nuevas coordenadas

$$z^{(j)} = (\mathbf{u}^{(j)})^\top (\mathbf{x} - \boldsymbol{\mu})$$

vemos que la expresión dentro de la exponencial toma la forma

$$-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) = \sum_{j=1}^D \frac{(z^{(j)})^2}{\lambda^{(j)}}$$

De aquí vemos fácilmente que las curvas de nivel de la densidad son **elipses**, ver Fig. 1.

2. Regresión lineal bayesiana

2.1. El modelo

Consideremos el vector (columna) de atributos $\mathbf{x} = [1, x^{(1)}, \dots, x^{(D)}]^\top$ y el vector (columna) de coeficientes $\boldsymbol{\theta} = [b, w_1, \dots, w_D]^\top$, de modo tal que nuestro modelo lineal se escribe

$$y = b + w_1 x^{(1)} + w_2 x^{(2)} + \dots + w_D x^{(D)} + \epsilon,$$

en donde $\epsilon \sim \mathcal{N}(0, \sigma^2)$. Recordar que en notación matricial esto es

$$y = \mathbf{x}^\top \boldsymbol{\theta} + \epsilon$$

Apilando todas las observaciones podemos calcular todas las predicciones:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta} + \boldsymbol{\epsilon}$$

en donde recordar que

$$\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^\top$$

es la matriz de diseño;

$$\mathbf{y} = [y_1, \dots, y_N]^\top$$

es el vector de targets o etiquetas, y

$$\boldsymbol{\epsilon} = [\epsilon_1, \dots, \epsilon_N]^\top$$

el vector de ruidos.

Vamos a sumir (solo para simplificar la exposición) que σ^2 es un parámetro en el sentido clásico, y que además conocemos su valor.

La regla Bayes:

$$p(\boldsymbol{\theta} | \mathbf{y}, \mathbf{X}) = \frac{p(\mathbf{y} | \boldsymbol{\theta}, \mathbf{X}) p(\boldsymbol{\theta})}{p(\mathbf{y} | \mathbf{X})}$$

Si expandimos la verosimilitud marginal:

$$p(\boldsymbol{\theta} | \mathbf{y}, \mathbf{X}) = \frac{p(\mathbf{y} | \boldsymbol{\theta}, \mathbf{X}) p(\boldsymbol{\theta})}{\int p(\mathbf{y} | \boldsymbol{\theta}, \mathbf{X}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}}.$$

Estamos interesados en hacer predicciones, por lo que para un vector de atributos $\mathbf{x}_{\text{nuevo}}$ la densidad de y_{nuevo} está dada por:

$$p(y_{\text{nuevo}} | \mathbf{x}_{\text{nuevo}}, \mathbf{X}, \mathbf{y}) = \int p(y_{\text{nuevo}} | \mathbf{x}_{\text{nuevo}}, \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathbf{y}, \mathbf{X}) d\boldsymbol{\theta}$$

Una vez calculada la densidad podemos dar predicciones:

$$P(a < y_{\text{nuevo}} \leq b | \mathbf{x}_{\text{nuevo}}, \mathbf{X}, \mathbf{y}) = \int_a^b p(z | \mathbf{x}_{\text{nuevo}}, \mathbf{X}, \mathbf{y}) dz$$

2.2. Verosimilitud

La verosimilitud $p(\mathbf{y} \mid \boldsymbol{\theta}, \mathbf{X})$ viene dada por nuestro modelo

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta} + \boldsymbol{\epsilon}$$

donde $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_N)$. Es decir:

$$p(\mathbf{y} \mid \boldsymbol{\theta}, \mathbf{X}) = \mathcal{N}(\mathbf{X}\boldsymbol{\theta}, \sigma^2 \mathbf{I}_N),$$

una Gaussiana multivariada en dimensión N , con media $\mathbf{X}\boldsymbol{\theta}$ y matriz de covarianzas $\sigma^2 \mathbf{I}_N$.

2.3. A priori

Vamos a elegir una distribución a priori $p(\boldsymbol{\theta})$ que sea conjugada a la verosimilitud:

$$p(\boldsymbol{\theta}; \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) = \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0),$$

cuyos parámetros son la media $\boldsymbol{\mu}_0$ y la matriz de covarianzas $\boldsymbol{\Sigma}_0$.

2.4. A posteriori

Vamos a calcular la distribución a posteriori:

$$\begin{aligned} p(\boldsymbol{\theta} \mid \mathbf{y}, \mathbf{X}) &\propto p(\mathbf{y} \mid \boldsymbol{\theta}, \mathbf{X}) p(\boldsymbol{\theta}) \\ &= \frac{1}{(2\pi)^{N/2} |\sigma^2 \mathbf{I}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\sigma^2 \mathbf{I})^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})\right) \\ &\quad \times \\ &\quad \frac{1}{(2\pi)^{N/2} |\boldsymbol{\Sigma}_0|^{1/2}} \exp\left(-\frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\mu}_0)^\top \boldsymbol{\Sigma}_0^{-1} (\boldsymbol{\theta} - \boldsymbol{\mu}_0)\right) \end{aligned}$$

Eliminando las constantes nos quedamos con:

$$\exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})\right) \exp\left(-\frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\mu}_0)^\top \boldsymbol{\Sigma}_0^{-1} (\boldsymbol{\theta} - \boldsymbol{\mu}_0)\right)$$

O de forma equivalente:

$$\exp\left\{-\frac{1}{2}\left(\frac{1}{\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) + (\boldsymbol{\theta} - \boldsymbol{\mu}_0)^\top \boldsymbol{\Sigma}_0^{-1} (\boldsymbol{\theta} - \boldsymbol{\mu}_0)\right)\right\}$$

Desarrollando y descartando todo lo que no tenga $\boldsymbol{\theta}$:

$$p(\boldsymbol{\theta} \mid \mathbf{y}, \mathbf{X}) \propto \exp\left\{-\frac{1}{2}\left(-\frac{2}{\sigma^2}\mathbf{y}^\top \mathbf{X}\boldsymbol{\theta} + \frac{1}{\sigma^2}\boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\theta} + \boldsymbol{\theta}^\top \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\theta} - 2\boldsymbol{\mu}_0^\top \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\theta}\right)\right\}$$

Sabemos que la distribución a posteriori es Gaussiana, digamos de parámetros $\boldsymbol{\mu}_1$ y $\boldsymbol{\Sigma}_1$:

$$\begin{aligned} p(\boldsymbol{\theta} \mid \mathbf{y}, \mathbf{X}) &= \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \\ &\propto \exp\left(-\frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\mu}_1)^\top \boldsymbol{\Sigma}_1^{-1}(\boldsymbol{\theta} - \boldsymbol{\mu}_1)\right) \\ &\propto \exp\left\{-\frac{1}{2}(\boldsymbol{\theta}^\top \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\theta} - 2\boldsymbol{\mu}_1^\top \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\theta})\right\} \end{aligned}$$

Igualando los términos cuadráticos podemos despejar $\boldsymbol{\Sigma}_1$:

$$\begin{aligned} \boldsymbol{\theta}^\top \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\theta} &= \frac{1}{\sigma^2} \boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\theta} + \boldsymbol{\theta}^\top \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\theta} \\ &= \boldsymbol{\theta}^\top \left(\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X} + \boldsymbol{\Sigma}_0^{-1} \right) \boldsymbol{\theta} \end{aligned}$$

de donde obtenemos:

$$\boldsymbol{\Sigma}_1 = \left(\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X} + \boldsymbol{\Sigma}_0^{-1} \right)^{-1}$$

De forma similar, igualando los términos lineales y usando la expresión para $\boldsymbol{\Sigma}_1$ obtenemos una expresión para $\boldsymbol{\mu}_1$:

$$\begin{aligned} -2\boldsymbol{\mu}_1^\top \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\theta} &= -\frac{2}{\sigma^2} \mathbf{y}^\top \mathbf{X} \boldsymbol{\theta} - 2\boldsymbol{\mu}_0^\top \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\theta} \\ \boldsymbol{\mu}_1^\top \boldsymbol{\Sigma}_1^{-1} &= \frac{1}{\sigma^2} \mathbf{y}^\top \mathbf{X} + \boldsymbol{\mu}_0^\top \boldsymbol{\Sigma}_0^{-1} \\ \boldsymbol{\mu}_1^\top &= \left(\frac{1}{\sigma^2} \mathbf{y}^\top \mathbf{X} + \boldsymbol{\mu}_0^\top \boldsymbol{\Sigma}_0^{-1} \right) \boldsymbol{\Sigma}_1 \end{aligned}$$

de donde obtenemos

$$\boldsymbol{\mu}_1 = \boldsymbol{\Sigma}_1 \left(\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{y} + \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0 \right)$$

En resumen:

$$p(\boldsymbol{\theta} \mid \mathbf{y}, \mathbf{X}) = \mathcal{N}\left(\left(\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X} + \boldsymbol{\Sigma}_0^{-1} \right)^{-1} \left(\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{y} + \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0 \right), \left(\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X} + \boldsymbol{\Sigma}_0^{-1} \right)^{-1} \right)$$

2.5. Comparación con la regularización Ridge

Si elegimos $\boldsymbol{\mu}_0 = \mathbf{0}$ y $\boldsymbol{\Sigma}_0 = \frac{\sigma^2}{N\alpha}\mathbf{I}$ obtenemos que la media de la distribución a posteriori es

$$\boldsymbol{\mu}_1 = (\mathbf{X}^\top \mathbf{X} + N\alpha\mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

que es igual a la solución del problema de regresión de Ridge:

$$\hat{\boldsymbol{\theta}}_\alpha = \begin{bmatrix} \hat{b} \\ \hat{\mathbf{w}} \end{bmatrix} = \arg \min_{\boldsymbol{\theta}} \left\{ \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\| + \alpha\|\mathbf{w}\|^2 \right\}.$$

Es decir, la regresión bayesiana ya tiene incluida una forma de regularización con la consideración de una distribución a priori.

2.6. Predicciones

Dada una observación nueva $\mathbf{x}_{\text{nuevo}}$, nos interesa la densidad:

$$p(y_{\text{nuevo}} \mid \mathbf{x}_{\text{nuevo}}, \mathbf{X}, \mathbf{y}).$$

Una forma de hacerlo es integrar respecto de $\boldsymbol{\theta}$ ponderando por la distribución a posteriori $p(\boldsymbol{\theta} \mid \mathbf{y}, \mathbf{X})$:

$$\begin{aligned} p(y_{\text{nuevo}} \mid \mathbf{x}_{\text{nuevo}}, \mathbf{X}, \mathbf{y}) &= \mathbb{E}_{\text{A posteriori}} \{p(y_{\text{nuevo}} \mid \mathbf{x}_{\text{nuevo}}, \boldsymbol{\theta})\} \\ &= \int p(y_{\text{nuevo}} \mid \mathbf{x}_{\text{nuevo}}, \boldsymbol{\theta}) p(\boldsymbol{\theta} \mid \mathbf{y}, \mathbf{X}) d\boldsymbol{\theta} \end{aligned}$$

Sin embargo, nuestro modelo es:

$$y_{\text{nuevo}} = \mathbf{x}_{\text{nuevo}}^\top \boldsymbol{\theta} + \epsilon$$

con $\boldsymbol{\theta} \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ y $\epsilon \sim \mathcal{N}(0, \sigma^2)$. Además estos términos son independientes.

Cómo la proyección $\mathbf{x}^\top \boldsymbol{\theta}$ es normal, la distribución de y_{nuevo} es también normal.

Su media es:

$$\mathbb{E}_{\text{A posteriori}} [y_{\text{nuevo}}] = \mathbb{E}_{\text{A posteriori}} [\mathbf{x}_{\text{nuevo}}^\top \boldsymbol{\theta}] = \mathbf{x}_{\text{nuevo}}^\top \mathbb{E}_{\text{A posteriori}} [\boldsymbol{\theta}] = \mathbf{x}_{\text{nuevo}}^\top \boldsymbol{\mu}_1$$

Su varianza es:

$$\mathbb{E}_{\text{A posteriori}} [y_{\text{nuevo}}^2] - (\mathbf{x}_{\text{nuevo}}^\top \boldsymbol{\mu}_1)^2 = \mathbb{E}_{\text{A posteriori}} [(\mathbf{x}_{\text{nuevo}}^\top \boldsymbol{\theta} + \epsilon)^2] - (\mathbf{x}_{\text{nuevo}}^\top \boldsymbol{\mu}_1)^2$$

El primer término del lado derecho es igual a

$$\begin{aligned}
 \mathbb{E}_{\text{A posteriori}} \left[(\mathbf{x}_{\text{nuevo}}^\top \boldsymbol{\theta})^2 \right] + \sigma^2 &= \mathbb{E}_{\text{A posteriori}} \left[\mathbf{x}_{\text{nuevo}}^\top \boldsymbol{\theta} \boldsymbol{\theta}^\top \mathbf{x}_{\text{nuevo}} \right] + \sigma^2 \\
 &= \mathbf{x}_{\text{nuevo}}^\top \mathbb{E}_{\text{A posteriori}} \left[\boldsymbol{\theta} \boldsymbol{\theta}^\top \right] \mathbf{x}_{\text{nuevo}} + \sigma^2 \\
 &= \mathbf{x}_{\text{nuevo}}^\top (\boldsymbol{\Sigma}_1 + \boldsymbol{\mu}_1 \boldsymbol{\mu}_1^\top) \mathbf{x}_{\text{nuevo}} + \sigma^2 \\
 &= \mathbf{x}_{\text{nuevo}}^\top \boldsymbol{\Sigma}_1 \mathbf{x}_{\text{nuevo}} + (\mathbf{x}_{\text{nuevo}}^\top \boldsymbol{\mu}_1)^2 + \sigma^2
 \end{aligned}$$

y finalmente obtenemos que la varianza es:

$$\mathbf{x}_{\text{nuevo}}^\top \boldsymbol{\Sigma}_1 \mathbf{x}_{\text{nuevo}} + \sigma^2$$

En resumen

$$p(y_{\text{nuevo}} \mid \mathbf{x}_{\text{nuevo}}, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{x}_{\text{nuevo}}^\top \boldsymbol{\mu}_1, \mathbf{x}_{\text{nuevo}}^\top \boldsymbol{\Sigma}_1 \mathbf{x}_{\text{nuevo}} + \sigma^2)$$

2.7. Verosimilitud marginal

La verosimilitud marginal puede utilizarse para seleccionar hiperparámetros. A modo de ejemplo comentaremos sobre su utilización en la selección del grado en una regresión polinomial.

Para un polinomio de grado K el modelo de ruido Gaussiano es:

$$y = b + w_1 x_n + w_2 x^2 + \dots + w_K x^K + \epsilon$$

donde $\epsilon \sim \mathcal{N}(0, \sigma^2)$.

En forma vectorial:

$$\mathbf{y} = \mathbf{X} \boldsymbol{\theta} + \boldsymbol{\epsilon},$$

donde $\boldsymbol{\theta} = [b, w_1, \dots, w_K]^\top$, $\mathbf{x} = [1, x, x^2, \dots, x^K]^\top$, $\mathbf{y} = [y_1, \dots, y_N]^\top$ y la matriz de diseño, $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^\top$.

La verosimilitud marginal está dada por la integral

$$p(\mathbf{y} \mid \mathbf{X}) = \int p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}.$$

De forma análoga a las cuentas que hemos hecho hasta ahora se puede ver que:

$$p(\mathbf{y} \mid \mathbf{X}) = \mathcal{N}(\mathbf{X} \boldsymbol{\mu}_0, \sigma^2 \mathbf{I}_N + \mathbf{X} \boldsymbol{\Sigma}_0 \mathbf{X}^\top)$$

Luego podemos seleccionar el modelo con mayor verosimilitud marginal.